

[Patent Document]

1. Japanese Patent Laid Open  
No. 2003-303055

Disk storage system having disk arrays connected with disk adapters through switches

Hitachi, Ltd.

Inventor(s): Tanaka, Katsuya ; Fujimoto, Kazuhisa

Application No. 10/212882, Filed 20020807, A1 Published 20031009

Abstract:

A disk storage system has high throughput between a disk adapter of a disk controller and a disk array. The disk adapter of the disk controller is connected to the disk array through switches. Data on a channel between the switch and a RAID group is multiplexed in the switch to be transferred onto a channel between the switch and the disk adapter and data on the channel between the switch and the disk adapter is demultiplexed in the switch to be transferred onto the channel between the switch and the RAID group. A data transfer rate on the channel between the disk adapter and the switch is made higher than that on the channel.

US.Class: 711114 711154



りデータ転送速度を第1のスイッチと前記複数のディスクアレイ間のチャネル当りデータ転送速度より高く設定し、

第2のディスクアダプタと第2のスイッチ間、および第1のディスクアダプタと第2のスイッチ間のチャネル当りデータ転送速度を第2のスイッチと前記複数のディスクアレイ間のチャネル当りデータ転送速度より高く設定し、

第1のスイッチと第2のスイッチを、第1のディスクアダプタと第2のスイッチ間を接続したチャネルと同等のデータ転送速度を有するチャネルと、第2のディスクアダプタと第1のスイッチ間を接続したチャネルと同等のデータ転送速度を有するチャネルと、を介して接続し、第1のスイッチは、第1のディスクアダプタまたは第2のディスクアダプタまたは第2のスイッチが接続されたポートと前記複数のディスクアレイが接続された各ポートとの間のポート間の接続の切り換えを、入力されたフレーム毎に、該フレーム内の送信用情報にしたがって行う、

第2のスイッチは、第1のディスクアダプタまたは第2のディスクアダプタまたは第1のスイッチが接続されたポートと前記複数のディスクアレイが接続された各ポートとの間のポート間の接続の切り換えを、入力されたフレーム毎に、該フレーム内の送信用情報にしたがって行うことを特徴とするディスク装置、

【請求項6】 請求項1乃至請求項5のいずれかの請求項記載のディスク装置において、

前記ディスクアレイからデータ読み出し時には、前記ディスクアダプタから前記スイッチに転送されるデータを前記スイッチにおいて逆多量化して前記ディスクアレイに転送することを特徴とするディスク装置、

【請求項7】 請求項1乃至請求項5のいずれかの請求項記載のディスク装置において、

ディスクアダプタからディスクアレイへのデータ書き込み時に、前記ディスクアダプタは、ラウンドロビン方式により前記ポート間の接続を切り替えることを特徴とするディスク装置、

【請求項8】 請求項7記載のディスク装置において、周期的に切り換えられるポート数を、ディスクアダプタとスイッチ間のチャネル当りデータ転送速度の、スイッチとディスクアレイ間のチャネル当りデータ転送速度に対する比、と同程度に設定することを特徴とするディスク装置

図、

【請求項9】 請求項1乃至請求項5のいずれかの請求項記載のディスク装置において、

前記ディスクアダプタと前記スイッチ間を光ファイバケーブルで接続し、前記スイッチと前記ディスクアレイ間をメタルケーブルで接続することを特徴とするディスク装置、

【発明の詳細な説明】

【0001】

【発明の属する技術分野】 本発明は、コンピュータシステムにおける2次記憶装置に関し、特に出入口データ転送性能が高いディスク装置に関する。

【0002】

【従来の技術】 現在のコンピュータシステムにおいては、CPU（中央処理装置）が必要とするデータは2次記憶装置に保存され、CPUなどが必要とするときに必要に応じて2次記憶装置に対してデータの書き込みおよび読み出しを行う。この2次記憶装置としては、一般に不揮発性記憶媒体が使用され、代表的なものとして磁気ディスク装置や、光ディスクなどのディスク装置がある。近年高度情報化に伴い、コンピュータシステムにおいて、この種の2次記憶装置の高性能化が要求されている。

【0003】 図9に、従来のディスク装置のブロック図を示す。図9において、ディスク装置はディスクコントローラDKCとディスクアレイDAで構成される。ディスクコントローラDKCは、上位側CPU（図示せず）とディスク装置を接続するチャネルアダプタCHAと、ディスクアレイDAに対して読み書きするデータを一時保存するキャッシュメモリCMと、ディスクコントローラDKCとディスクアレイDAを接続するディスクアダプタDCAからなる。チャネルアダプタCHAとキャッシュメモリCMとディスクアダプタDCAは、バスまたはスイッチで相互接続されている。チャネルアダプタCHAはC1、C2、C3、C4の4本のチャネルでCPUと接続している。ディスクアダプタDCAはD1、D2、D3、D4の4本のチャネルでディスクアレイと接続している。ここでディスクアレイDAはディスクグループR1、R2、R3、R4からなり、それぞれRAIDグループを構成する。

【0004】 チャネルC1、C2、C3、C4から入力された書き込みデータは、キャッシュメモリCMに接続された書き込みと同時に、該データをブロックサイズ単位に分割し、チャネルD1、D2、D3、D4の4本のチャネルにはブロック単位に分割されたデータ、残りの1チャネルは前記分割データから計算したパリティを、ディスクアダプタDCAから計算したパリティを、ディスクアダプタDCAからディスクアレイDAへ送る。データ読み出し時は、まずキャッシュメモリCM内に該当データの有無を調べる。有る場合は、キャッシュメモリCMからディスクアレイDAへデータを送る。

【0005】 請求項1乃至請求項5のいずれかの請求項記載のディスク装置において、

ディスクアダプタからディスクアレイへのデータ書き込み時に、前記ディスクアダプタは、ラウンドロビン方式により前記ポート間の接続を切り替えることを特徴とするディスク装置、

【請求項6】 請求項1乃至請求項5のいずれかの請求項記載のディスク装置において、

ディスクアダプタからディスクアレイへのデータ書き込み時に、前記ディスクアダプタは、ラウンドロビン方式により前記ポート間の接続を切り替えることを特徴とするディスク装置、

キャッシュメモリCMからディスクアレイDAへデータを送る。第2の従来技術では、ディスクアダプタとディスクアレイとの間にスイッチを適用することによりディスク増設ポート数を増加させることができるが、チャネル当りのデータ転送速度はディスクアレイのデータ転送速度に制限されるので、ディスクアダプタとディスクアレイ間のスループットが性能ネックになるという問題があった。第3の従来技術は、ディスクの回転待ち時間の影響を低減できる技術であり、フロントエンドとバックエンドのスループット増強は低減できないという問題があった。

【0006】 ディスクアダプタとディスクアレイを、スイッチを介して接続したディスク装置が、特開平5-173722号の「マルチチャネルデータおよびパリティの交換デバイス」に開示されている。以下、該公開に記載の従来技術を第9の従来技術と呼ぶ。第2の従来技術によれば、ディスクアレイに開通したバス本数とディスクアダプタに開通したバス本数とを独立に設定できる。ディスクアダプタとディスクアレイを、バックエンドの接続を介して接続したディスク装置が、特開平6-196277号の「回転遅延装置」に開示されている。以下、該公開に記載の従来技術を第3の従来技術と呼ぶ。第3の従来技術によれば、ディスクアダプタとディスクアレイ間のデータ転送速度を任意に設定でき、ディスクの回転待ちの影響を低減できる。

【0007】 請求項1乃至請求項5のいずれかの請求項記載のディスク装置において、

ディスクアダプタからディスクアレイへのデータ書き込み時に、前記ディスクアダプタは、ラウンドロビン方式により前記ポート間の接続を切り替えることを特徴とするディスク装置、

【請求項6】 請求項1乃至請求項5のいずれかの請求項記載のディスク装置において、

ディスクアダプタからディスクアレイへのデータ書き込み時に、前記ディスクアダプタは、ラウンドロビン方式により前記ポート間の接続を切り替えることを特徴とするディスク装置、

スイッチは、ディスクアダプタが接続されたポートと同一の数のディスクアレイが接続された各ポートとの間でのポート間の接続の切り換えを、入力されたフレーム毎に、該フレーム内の送信先情報にしたがって行っている。また、前記ディスクアダプタと前記ディスクアレイを、パッドメモリを有するスイッチを介して接続し、同一のスロットに接続したディスクドライブの組み合わせでRAIDグループを構成し、ディスクアダプタとスロット間のチャネル当りデータ転送速度を、スイッチとディスクアレイ間のチャネル当りデータ転送速度より高く設定し、スイッチは、ディスクアダプタが接続されたポートとRAIDグループを構成するディスクドライブが接続された各ポートとの間でポート間の接続の切り換えを、入力されたフレーム毎に、該フレーム内の送信先情報にしたがって行っている。また、第1のディスクコントローラと第2のディスクコントローラと複数のディスクアレイからなり、第1のディスクメモリと第1のチャネルアダプタと第1のキャッシュメモリと第1のディスクアダプタを有し、第2のディスクコントローラは第2のチャネルアダプタと第2のキャッシュメモリと第2のディスクアダプタを有するディスク装置であり、第1のディスクアダプタと前記複数のディスクアレイとをバッファメモリを有する第1のスイッチを介して接続し、且つ第2のディスクアダプタと前記複数のディスクアレイとをバッファメモリを有する第2のスイッチを紹介して接続し、さらに第1のスイッチと第2のディスクアダプタを接続し、第2のスイッチと第1のディスクアダプタを接続し、第2のディスクアダプタと第2のスイッチと前記第1のディスクアダプタと第2のディスクアダプタとの間のチャネル当りデータ転送速度を第2のスイッチと前記複数のディスクアレイ間のチャネル当りデータ転送速度より高く設定し、第1のスイッチは、第1のディスクアダプタまたは第2のディスクアダプタが接続されたポートと前記複数のディスクアダプタが接続された各ポートとの間でポート間の接続の切り換えを、入力されたフレーム毎に、該フレーム内の送信先情報にしたがって行っている。また、上記第1のディスクアダプタと第2のスイッチ間を接続したチャネル同等のデータ転送速度を有するチャネルと、第2のディスクアダプタと第1のスイッチ間を接続したチャネル同等のデータ転送速度を有するチャネルとを、介して接続している。また、前記ディスクアレイからのデータ読み出し時には、前記ディスクアレイから前記第1のディスクアダプタに転送されるデータを前記スイッチにおいて多重化して前記ディスクアダプタに転送し、前記ディ

スクアレレイへのデータ書き込み時には、前記ディスタクタクタから前記スイッチに転送されるデータを前記スイッチにおいて逆多重化して前記ディスタクアレイに転送するようにしている。また、前記ディスタクタクタから前記ディスタクアレイへのデータ書き込み時に、前記ディスタクタは、前記ポート間の接続の切り替えが行われるように、送出するフレームに送信先情報を設定し、前記ディスタクアレイから前記ディスタクタへのデータを読み出し時に、前記スイッチは、ラウンドロビン方式により前記ポート間の接続を切り替えるようにしている。また、さらに、切り替えるポート数を、ディスタクタとスイッチ間のチャネル当りデータ転送速度の、スイッチとディスタクアレイ間のチャネル当りデータ転送速度に対する比、と同程度に設定している。また、前記ディスタクタと前記スイッチ間を光ファイバケーブルで接続し、前記スイッチと前記ディスタクアレイ間をメタルケーブルで接続するようにしている。

[0010]

【発明の実施の形態】以下、図面を参照して本発明の実施の形態を詳細に説明する。図1に本発明の、第1の実施の形態であるディスタクタ装置の構成を示す。本実施の形態のディスタクタ装置は、ディスタクトローラDKCとディスタクアレイDAからなる。ディスタクトローラDKCは、チャネルアダプタCHAと、キャプシユメモリCMと、ディスタクタダブダDKAからなる。チャネルアダプタCHAは、上位CPU（図示せず）とディスタクトローラDKCとがデータを送受する際の制御を行う。C1、C2、C3およびC4は、チャネルアダプタCHAがCPUと通信するチャネルである。キャプシユメモリCMは、本実施の形態のディスタクタ装置が入力するデータを一時保存するメモリである。ディスタクタDKAは、ディスタクトローラDKCとディスタクアレイDAとがデータを送受する際の制御を行う。ディスタクタダブダDKAは、チャネルD01、D02、D03、D04を介して、ディスタクアレイDAと接続する。ディスタクタダブダDKAとディスタクアレイDAは、チャネルD01、D02、D03、D04上で全二重通信が可能である。

[0011] ここで、本実施の形態のディスタクタ装置は、ディスタクタダブダDKAとディスタクアレイDAを、スイッチSW1、SW2、SW3、SW4を介して接続している点に特徴がある。ディスタクアレイDAは、ディスタクトローラR1、R2、R3、R4からなる。ディスタクトローラR1は、スイッチSW1介してディスタクタダブダDKAと接続する。同様に、ディスタクトローラR2はスイッチSW2を介して、ディスタクトローラR3はスイッチSW3を介して、ディスタクトローラR4はスイッチSW4を介して、それぞれディスタクタダブダDKAと接続する。

[0012] 本実施の形態のディスタクタ装置においてRA

2、in3、in4、in5と、出力ポートout1、out2、out3、out4、out5を有する。ポートP1から入力されたフレームは、シリアルパラレル変換器SP12と、バッファメモリBM1と、8B10B変換デコーダDEC1を經由し、スイッチコントローラCTLと入力ポートin1へ入力される。スイッチコントローラCTLにおいて、入力フレームのヘッダ部分に格納された送信先アドレスを解釈し、クロスバスイッチXSWを切り換える。例として、ポートP2が出力先として選ばれた場合は、入力されたフレームは出力ポートout2と、8B10B変換エンコーダENC2と、バッファメモリBM2と、シリアルシリアル変換器SP2を經由して、ポートP2から出力される。ここで、バッファメモリBM1、BM2はFIFO (First-In-First-Out) メモリである。

{0016} シリアルパラレル変換器SP1は、8B10B符号化されたシリアルデータを10ビット幅のラレルデータに変換し、ポートP1におけるデータ転送速度の1/10の速度に同期してバッファメモリBM1に書き込む。8B10BデコーダDEC1は、クロスバスイッチXSWの動作速度に同期して、10ビットラレルデータをバッファメモリBM1から読み出し、8B10B復号化して、8ビットラレルデータに変換する。8B10BエンコーダENC2は、クロスバスイッチXSWでスイッチされた8ビットラレルデータを再び8B10B符号化し、10ビットラレルデータに変換後、クロスバスイッチXSWの動作速度に同期してバッファメモリBM2に書き込む。パラレルシリアル変換器SP2は、ポートP2におけるデータ転送速度の1/10の速度に同期して、10ビットラレルデータで、バッファメモリBM2から読み出し、シリアル化して、ポートP2から出力する。以上によりスイッチSW1は、ポートP1におけるデータ転送速度からポートP2におけるデータ転送速度へ速度変換する。

{0017} 図4は、ポートP1へ入力するフレームと、ポートP2、P3、P4、P5から出力されるフレームを示した図である。波形の凸はフレームが存在している時間、凹はフレームが存在していない時間を示している。フレームは伝送するデータ容量に従ってそのフレーム長が変化するが、ここではディスプレイレイヘのシグナルアクセスが行われており、フレーム長が一定である。図4では、入力ポートP1でのデータ転送速度が出力ポートP2、P3、P4、P5におけるデータ転送速度のm倍あるとする。従って、ポートP1におけるフレームFb2の時間T1は、ポートP2からの出力時間にT3へ伸びている。ここでT3=m×T1である。

{0018} 入力データのデータ転送速度が強く、且つ出力のデータ転送速度が遅い場合は、スイッチがスループットの低下無く換えないと出力ポートのバッファメモリが溢れ、スループットが低下する。フレームがスループットの低下無く

スイッチを通する。図4のように周期的に出力ポートを切り替える必要がある。スイッチ切り替えポート数をnとすると、スイッチ切り替え間隔T2はn×T1である(フレームの無い場合は無視した)。T2≧T3ならば、フレームの面無、スループットの低下は起こらない。T2≧T3はn≧mと同じである。つまり、ディスクアレイへのデータを送る時に、スイッチにおいてスループット低下を起こさないための条件は、周期的に切り替えるスイッチポート数nを、ディスクアダプタとスイッチ間のチャネル当りデータ転送速度の、スイッチとディスクアレイ間のチャネル当りデータ転送速度に対する比m、以上に設定することである。この条件が保たれるば、スイッチSW1は、ポートP1から入力したデータをバッファメモリにおいて速度変換し、フレーム単位で周期的に切り替えることにより逆多量化し、ポートP2、P3、P4、P5へ分配して出力する。スイッチを、周期的に切り替える方法の一つは、スイッチに接続した、ディスクグループをRAIDグループとすることである。RAIDのストライピング原理に従えば、スイッチは周期的に切り替わる。

[0019] 図3は、ポートP2、P3、P4、P5からフレームを入力し、ポートP1から出力する場合を示す。これはディスクアレイからの読み出し時のスイッチ動作に相当する。例えば、ポートP2から入力したフレームは、シリアルパラレル変換回路SP2から、バッファメモリBM2と、8B10B変換デコードDEC2を経由し、スイッチコントローラCTLと入力ポートIn2へ入力される。スイッチコントローラCTLにおいて、入力フレームのヘッダ部分に書かれた送信先アドレスを解読し、クロスバスイッチXSWを切り替える。図3の場合、ラウンドロビン方式によりクロスバスイッチXSWを切り替えて、順番にポートP2、P3、P4、P5から入力されるデータは全てポートP1へ出力する。すなわち、読み出し時は、複数の入力ポート(P2、P3、P4、P5)に同時にフレームが流れる。これら複数の入力フレームは同時に入力ポートに届く必要はない。スイッチは、総当り的に入力ポート間接続を切り替えることにより、これら複数の入力フレームを1フレームずつ出力ポート(P1)へ転送する。このように、スイッチを総当り的に切り替える方式を、ラウンドロビン(Round Robin)方式と呼ぶ。ラウンドロビン方式により、結果的にスイッチは周期的に切り替わることになる。なお、読み出し時においても、スイッチはフレーム内送信先情報に従って切り替わることに違いない。フレームは出力ポートOut1と、8B10B変換エンコードDEC1と、バッファメモリBM1と、パラレルシリアル変換回路SP1を経由して、ポートP1から出力される。

[0020] シリアルパラレル変換回路SP2は、8B10B符号化されたシリアルデータを10b11幅のバ

ラレルデータに変換し、ポートP2におけるデータ転送速度の1/10の速度に同期してバッファメモリBM2に書き込む。8B10BデコードDEC2は、クロスバスイッチXSWの動作速度に同期して、10b11パラレルデータをバッファメモリBM2から読み出し、8B10B復号化して、8b11パラレルデータに変換する。8B10BエンコードENC1は、クロスバスイッチXSWでスイッチされた8b11パラレルデータを再び8B10B符号化し、10b11パラレルデータに変換後、クロスバスイッチXSWの動作速度に同期してバッファメモリBM1に書き込む。パラレルシリアル変換回路PS1は、ポートP1におけるデータ転送速度の1/10の速度に同期して、10b11パラレルデータをバッファメモリBM1から読み出し、シリアル化して、ポートP1から出力する。以上によりスイッチSW1は、ポートP2におけるデータ転送速度からポートP1におけるデータ転送速度へ速度変換する。

[0021] 図5は、ポートP2、P3、P4、P5へ入力するフレームと、ポートP1から出力されるフレームを示した図である。波形の凸はフレームが存在する時間、凹はフレームが存在していない時間を示している。フレームは伝送するデータ容量に従ってそのフレーム長が変化するが、ここではディスクアレイへのシーケンシャルアクセスが行われており、フレーム長が一定である。図5では、入力ポートP1でのデータ転送速度が出力ポートP2、P3、P4、P5におけるデータ転送速度のm倍あるとする。従って、ポートP5におけるフレームFe5の時間T4は、ポートP1からの出力時にT5へ縮んでいる。ここでT4=m×T5である。フレームFe2、Fe3、Fe4、Fe5をポートP1から出力するのにかかる時間をT6とする。スイッチ切り替えポート数をnとすると、T6=n×T5である(フレームの無い場合は無視した)。スイッチにおいて幅幅によるスループット低下を防止するために、T6≦T4とする必要がある。T6≦T4はn≦mと同じである。

[0022] つまり、ディスクアレイからのデータ読み出し時に、スイッチにおいてスループット低下を起こさないための条件は、周期的に切り替えるスイッチポート数nを、ディスクアダプタとスイッチ間のチャネル当りデータ転送速度の、スイッチとディスクアレイ間のチャネル当りデータ転送速度に対する比m、以下に設定することである。この条件が保たれるば、スイッチSW1は、ポートP2、P3、P4、P5から入力したデータをバッファメモリにおいて速度変換し、フレーム単位で周期的に切り替えることにより逆多量化し、ポートP1へ出力する。よって、ディスクアレイへの書き込みおよびディスクアレイからの読み出しを高スループット化するために、n≦m、つまり、周期的に切り替えるポート数を、ディスクアダプタとスイッチ間のチャネル当りデータ転送速度の、スイッチとディスクアレイ間のチャ

ネル当りデータ転送速度に対する比、と同程度に設定すればよいことが分かる。

[0023] 例えば、ディスクアダプタとスイッチ間の4Gbpsのチャネル1本で接続し、スイッチとディスクアレイ間を1Gbpsのチャネル4本で接続する。また、ディスクアダプタとスイッチ間の10Gbpsのチャネル1本で接続し、スイッチとディスクアレイ間を2Gbpsのチャネル4本で接続する。この場合、スイッチ入出力ポート間でスループットのバランスが取れないので、実効的なスループットは2Gbps×4=8Gbpsとなる。

[0024] 以上より、スイッチSW1において速度変換と逆多量化、逆多量化が行われるので、チャネルD1、D12、D13、D14上のデータ転送速度が低速でも、チャネルD01、D02、D03、D04でのデータ転送速度は高速にできる。つまり、ディスクアダプタDKAとディスクアレイDA間のスループットを上向きでできる。本実施の形態のディスク装置におけるデータ転送方式としては、ファイバチャネルやインフィニバンドが使用できる。

[0025] 図6は、第1の実施の形態のディスク装置において、ディスクドライブの増設方法を示した図である。図6では図1に対して、ディスクグループR5とR6が増設されている。ディスクドライブを増設するだけ、スイッチSW1とSW2としてポート数の多いスイッチを使用している。ディスクドライブを増設すると、スイッチのディスクアレイ側のスループットが増加し、ディスクアダプタ側のスループットバランスが崩れるので、スイッチの速度変換機能が有効に働かなくなる可能性がある。そこでスイッチSW1では、ディスクアダプタDKAとの間に、新規チャネルD05を増設している。また、スイッチSW2の場合は新規チャネルを増設せず、チャネルD02の信号伝送速度を増加させることで、ディスクアダプタ側とディスクアレイ側のスループットバランスを取っている。例えばスイッチSW1では、スイッチとディスクアレイ間を1Gbpsのチャネル8本で接続し、ディスクアダプタとスイッチ間を4Gbpsのチャネル2本で接続する。スイッチSW2では、スイッチとディスクアレイ間を1Gbpsのチャネル8本で接続し、ディスクアダプタとスイッチ間を10Gbpsのチャネル1本で接続する。このように、本実施の形態のディスク装置は、スイッチのポート数に応じたドライブ増設方法は、1ポート当たり接続できるドライブ数が少ないATA(AT Attachment)方式ディスクドライブを増設するのに適している。

[0026] 図7に本発明の、第2の実施の形態であるディスク装置の構成を示す。本実施の形態のディスク装置は、第1の実施の形態のディスク装置に対して、ディスクアレイ部分の構成方法が異なる。本実施の形態のデ

ィスク装置は、ディスクコントローラDKCと、4個のディスクアレイDA1、DA2、DA3、DA4からなる。ディスクコントローラDKCは、チャネルアダプタCHA、キャッシュメモリCM、ディスクアダプタDKAからなる。ディスクアレイDA1とディスクアダプタDKAは、チャネルD01とスイッチSW1を介して接続する。同様に、ディスクアレイDA2はチャネルD02とスイッチSW2を介して、ディスクアレイDA3はチャネルD03とスイッチSW3を介して、ディスクアレイDA4はチャネルD04とスイッチSW4を介して、それぞれディスクアダプタDKAと接続する。スイッチSW1、SW2、SW3とSW4は、第1の実施の形態と同様に速度変換と逆多量化、逆多量化を行うスイッチとして機能する。本実施の形態におけるディスクアダプタDKAと、スイッチSW1、SW2、SW3、SW4はファイバチャネルスイッチである。

[0027] 本実施の形態におけるディスクアレイの構成を、ディスクアレイDA1を例に述べる。ディスクアレイDA1、DA2、DA3、DA4は、同様のドライブ構成である。ディスクアレイDA1は、チャネルD11上に接続した4個のディスクからなるディスクアレイと、D12上に接続した4個のディスクからなるディスクアレイと、D13上に接続した4個のディスクからなるディスクアレイと、D14上に接続した4個のディスクからなるディスクアレイと、からなる。チャネルD11を例にとると、ディスクドライブDK1、DK2、DK3、DK4が、チャネルD11上に接続されている。このように、多数のドライブを一つのチャネル上に接続してディスクドライブにアクセスする方法としては、ファイバチャネルアービトラリティグループ(以下FC-A-Lと呼ぶ)がある。

[0028] 図10に、FC-A-Lの接続形態をディスクドライブDK1、DK2、DK3、DK4の接続形態を例として示す。各ディスクドライブの入出力ポートおよびスイッチSW1の入出力ポートは、送信機Txと受信機Rxを有する。FC-A-Lの接続形態は、例えば図10に示すように、各ドライブの入出力ポートおよびスイッチの入出力ポートをループ状に接続するポートロブである。各ドライブの入出力ポートはファイバチャネル(Node Loop)ポートとして機能する。NLポートとは、ループ動作をする装置(ここではディスクアレイDA1接続側)出力ポートは、ファイバチャネルのFL(Fabric Loop)ポートとして機能する。FLポートとは、FC-A-Lを接続可能なスイッチのポートである。FLポートを有するループは、ファイバチャネルのパブリックループとして機能するので、



チャネルD11が形成するFC-Aはパブリックルー  
プとなる。パブリックループとは、ループ上のディスク  
ドライバが、スイッチを介してループ外のポートと通信  
可能なループである。よって、ディスクドライバDK  
1、DK2、DK3、DK4は、スイッチSW1および  
チャネルD01を介してディスクアダプタDKA1と通信  
可能である。以上、チャネルD11の接続形態を例に取  
明したが、チャネルD12、D13、D14でも同様で  
ある。本実施の形態のディスク装置においてRAIDシ  
ステムを構成する場合は、ディスクグループR1、R  
2、R3、R4を、それぞれRAIDグループとす。  
本実施の形態では、4個のディスクドライバでRAID  
グループを構成しているが、RAIDグループを構成す  
るドライブ数を4個に限るものではない。

[0029] 本実施の形態においては、チャネルD1  
1、D12、D13、D14において、それぞれFC-  
Aを用いてディスクドライバを接続している。FC-  
Aの仕様に、チャネルD11、D12、D13、D  
14上には、それぞれ最大126台までのディスクドラ  
イバが接続可能である。また、チャネルD01、D0  
2、D03、D04の媒体として光ファイバケーブル  
を、チャネルD11、D12、D13、D14の媒体と  
してメタルケーブルを用いる。

[0030] 以上説明するように、本実施の形態のディ  
スク装置においては、ディスクドライバをFC-Aで  
接続しているため、スイッチのポート当りに接続できる  
ドライブ台数が増加できる。つまり、ディスク装置当り  
の記憶容量を増加させる効果がある。また、ディスクド  
ライバをメタルケーブルで接続することにより、ディス  
クドライブ毎に高価な光インターフェースを装備する必  
要がなくなるので、ディスクドライバのコストを下げ  
効果がある。

[0031] 図8に本発明の、第3の実施の形態である  
ディスク装置の構成を示す。本実施の形態のディスク装  
置は、ディスクコントローラとスイッチを二重化した点  
に特徴がある。本実施の形態において、ディスクアダ  
プタDKA1、DKA2と、スイッチSW1、SW2と、  
ディスクアレイDA1との間のデータ転送方式は、フ  
ァイバチャネルを使用している。本実施の形態のディ  
スク装置は、ディスクコントローラDKC1、DKC2と、  
スイッチSW1、SW2と、ディスクアレイDA1から  
なる。スイッチSW1とSW2は、第1の実施の形態と  
同様に速度変換と多重化、逆多重化を行うスイッチとし  
て機能する。ディスクコントローラDKC1は、チャネ  
ルアダプタDKA1と、キャッシュメモリCM1と、デ  
ィスクアダプタDKA1からなる。ディスクコントロー  
ラDKC2は、チャネルアダプタDKA2と、キャッ  
シュメモリCM2と、ディスクアダプタDKA2からな  
る。ディスクアダプタDKA1とスイッチSW1をチャ  
ネルD1aで接続し、ディスクアダプタDKA2とスイ

ッチSW2をチャネルD2aで接続し、ディスクアダ  
プタDKA1とスイッチSW2をチャネルD1bで接続  
し、ディスクアダプタDKA2とスイッチSW1をチャ  
ネルD2bで接続する。

[0032] ディスクアレイDA1を構成するディスク  
ドライバは、入出力ポートを2個有する。例えば、ディ  
スクドライバDK1、DK2、DK3、DK4は、チャ  
ネルD11およびD21の両チャネルと接続する。ディ  
スクアレイDA1は、チャネルD11とD21に接続し  
た4個のディスクからなるディスクアレイと、D12と  
D22に接続した4個のディスクからなるディスクアレ  
イト、D13とD23に接続した4個のディスクからな  
るディスクアレイと、D14とD24に接続した4個の  
ディスクからなるディスクアレイ、からなる。チャネル  
D11、D12、D13、D14、D21、D22、D  
23、D24は、FC-Aでディスクドライブを接続  
する。

[0033] 図11に本実施の形態におけるFC-A  
の接続形態を、ディスクドライブDK1、DK2、DK  
3、DK4の接続形態を例として示す。各ディスクド  
ライバは、それぞれNLポートを2個有する。各ディス  
クドライブの入出力ポートおよびスイッチSW1、SW2  
の入出力ポートは、送信機Txと受信機Rxを有する。

[0034] ディスクアレイDA1内の全ディスクドラ  
イバは、ディスクアダプタDKA1およびDKA2のど  
ちらからでもアクセス可能である。本実施の形態のディ  
スク装置は、チャネルD1b、D2bをスイッチSW  
1、SW2故障時の迂回経路として使用する。例えばス  
イッチSW1が故障した場合でも、ディスクアダプタD  
KA1はチャネルD1bとスイッチSW2経由でディス  
クアレイDA1にアクセスできる。逆に、スイッチSW

2が故障した場合は、ディスクアダプタDKA2はチャ  
ネルD2bとスイッチSW1経由でディスクアレイDA  
1にアクセスできるので、信頼性が高いディスク装置が  
実現できる。

[0035] 図12に本発明の、第4の実施の形態であ  
るディスク装置の構成を示す。本実施の形態のディス  
ク装置は、第3の実施の形態のディスク装置に對して、ス  
イッチSW1、SW2間を接続するチャネルD3a、D  
3bを設けた点に特徴がある。本実施の形態において、  
ディスクアダプタDKA1、DKA2と、スイッチSW  
1、SW2と、ディスクアレイDA1との間のデータ転  
送方式は、ファイバチャネルを使用している。本実施の  
形態のディスク装置は、ディスクコントローラDKC  
1、DKC2と、スイッチSW1、SW2と、ディス  
クアレイDA1からなる。スイッチSW1とSW2は、第  
1の実施の形態と同様に速度変換と多重化、逆多重化を  
行うスイッチとして機能する。ディスクコントローラD  
KC1は、チャネルアダプタDKA1と、キャッシュメ  
モリCM1と、ディスクアダプタDKA1からなる。デ  
ィスクコントローラDKC2は、チャネルアダプタDK  
A2と、キャッシュメモリCM2と、ディスクアダプタ  
DKA2からなる。ディスクアダプタDKA1とスイ  
ッチSW1をチャネルD1aで接続し、ディスクアダプタ  
DKA2とスイッチSW2をチャネルD2aで接続し、  
ディスクアダプタDKA1とスイッチSW2をチャネル  
D1bで接続し、ディスクアダプタDKA2とスイッチ  
SW1をチャネルD2bで接続する。さらに、スイ  
ッチSW1とSW2をチャネルD3a、D3bで接続する。  
[0036] ディスクアレイDA1を構成するディス  
クドライブは、入出力ポートを2個有する。例えば、ディ  
スクドライブDK1、DK2、DK3、DK4は、チャ  
ネルD11およびD21の両チャネルと接続する。ディ  
スクアレイDA1は、チャネルD11とD21に接続し  
た4個のディスクからなるディスクアレイと、D12と  
D22に接続した4個のディスクからなるディスクアレ  
イト、D13とD23に接続した4個のディスクからな  
るディスクアレイと、D14とD24に接続した4個の  
ディスクからなるディスクアレイ、からなる。チャネル  
D11、D12、D13、D14、D21、D22、D  
23、D24は、図11に示すようにFC-Aでディ  
スクドライブを接続する。ディスクアレイDA1内の全  
ディスクドライブは、ディスクアダプタDKA1および  
ディスクドライブは、ディスクアダプタDKA1および  
DKA2のどちらからでもアクセス可能である。本実施  
の形態のディスク装置においてRAIDシステムを構成  
する場合は、ディスクグループR1、R2、R3、R4  
を、それぞれRAIDグループとして使用する。本実施の形態で  
は、4個のディスクドライブでRAIDグループを構成  
しているが、RAIDグループを構成するドライブ数を  
4個に限るものではない。

[0037] ディスクアダプタDKA1、DKA2とデ

ィスクアレイDA1のアクセス経路について、先ず、定  
常時（スイッチ故障無しの場合）について説明する。デ  
ィスクアダプタDKA1は、チャネルD1aとSW1を  
介してディスクアレイDA1にアクセスする経路（経路  
05 1）と、チャネルD1bとスイッチSW2とチャネルD  
3aとスイッチSW1を介してディスクアレイDA1に  
アクセスする経路（経路2）を有する。同様に、ディ  
スクアダプタDKA2は、チャネルD2aとSW2を介し  
てディスクアレイDA1にアクセスする経路（経路3）  
と、チャネルD2bとスイッチSW1とチャネルD3b  
とスイッチSW2を介してディスクアレイDA1にアク  
セスする経路（経路4）を有する。一方、スイッチ故障  
時は、チャネルD1b、D2bを迂回経路として使用する。  
例えばスイッチSW1が故障した場合でも、ディス  
クアダプタDKA1はチャネルD1bとスイッチSW2  
経由でディスクアレイDA1にアクセスできる。逆に、  
スイッチSW2が故障した場合は、ディスクアダプタD  
KA2はチャネルD2bとスイッチSW1経由でディス  
クアレイDA1にアクセスできる。

[0038] 次に、本実施の形態におけるディスクアダ  
プタ-ディスクアレイ間のスループットについて説明す  
る。例として、チャネルD1a、D1b、D2a、D2  
b、D3a、D3b上のデータ伝送速度をチャネル当り  
2Gbpsとし、チャネルD11、D12、D13、D  
25 14、D21、D22、D23、D24上のデータ伝送  
速度をチャネル当り1Gbpsであるとする。このと  
き、スイッチSW1とディスクアレイDA1間の総スル  
ープットは4Gbpsである。ディスクアダプタDKA  
1とスイッチSW1間は、上記経路1および経路2でア  
クセスすることにより、総スループットは4Gbpsと  
なる。スイッチSW1のディスクアダプタDKA1側  
と、ディスクアレイDA1側のスループットが共に4G  
bpsであるので、ディスクアダプタDKA1とディス  
クアレイDA1間のスループットは4Gbpsとなる。

同様に、スイッチSW2とディスクアレイDA1間の総  
スループットは4Gbpsである。ディスクアダプタD  
KA2とスイッチSW2間は、上記経路3および経路4  
でアクセスすることにより、総スループットは4Gbp  
sとなる。スイッチSW2のディスクアダプタDKA2  
側と、ディスクアレイDA1側のスループットが共に4  
Gbpsであるので、ディスクアダプタDKA2とディ  
スクアレイDA2間のスループットは4Gbpsとなる。  
第3の実施の形態（図8）において、上記のチャネル当  
りスループット値を適用すると、チャネルD1b、D2  
bをスイッチ故障時の迂回経路として使用しないので、  
スループットは、チャネルD1a上のスループットに例  
えられ、2Gbpsとなる。同様に、ディスクアダプタ  
DKA2とディスクアレイDA1間のスループットは、  
チャネルD2a上のスループットに例えられ、2Gbps

となる。ディスクアダプタDKA1、DKA2とデ  
ィスクアレイDA1のアクセス経路について、先ず、定  
常時（スイッチ故障無しの場合）について説明する。デ  
ィスクアダプタDKA1は、チャネルD1aとSW1を  
介してディスクアレイDA1にアクセスする経路（経路  
05 1）と、チャネルD1bとスイッチSW2とチャネルD  
3aとスイッチSW1を介してディスクアレイDA1に  
アクセスする経路（経路2）を有する。同様に、ディ  
スクアダプタDKA2は、チャネルD2aとSW2を介し  
てディスクアレイDA1にアクセスする経路（経路3）  
と、チャネルD2bとスイッチSW1とチャネルD3b  
とスイッチSW2を介してディスクアレイDA1にアク  
セスする経路（経路4）を有する。一方、スイッチ故障  
時は、チャネルD1b、D2bを迂回経路として使用する。  
例えばスイッチSW1が故障した場合でも、ディス  
クアダプタDKA1はチャネルD1bとスイッチSW2  
経由でディスクアレイDA1にアクセスできる。逆に、  
スイッチSW2が故障した場合は、ディスクアダプタD  
KA2はチャネルD2bとスイッチSW1経由でディス  
クアレイDA1にアクセスできる。

8 となる。第3の実施の形態において、ディスクアダプタとディスクアレイ間スルーポートを4Gbpsにするためには、チャネルD1aおよびD2aのデータ伝送速度を、それぞれ4Gbpsに高める必要がある。以上から、本実施の形態によれば、ディスクアダプタとスイッチ間のチャネル当りデータ伝送速度が低くても、ディスクアダプタとディスクアレイ間の総スルーポートが高いディスク装置が実現できる。

【0039】

【発明の効果】 以上説明したように、本発明によれば以下の効果がある。ディスクアダプタとディスクアレイ間のスルーポートが高いディスク装置を提供できる。また、ディスクアダプタとディスクアレイ間のスルーポートが高く、且つディスクドライブ接続台数が多いディスク装置を提供できる。また、信頼性の高いディスクアレイを有するディスク装置を提供できる。また、信頼性が高いディスクアダプタとディスクアレイ間ネットワークを有するディスク装置を提供できる。また、信頼性およびスルーポートが高いディスクアダプタとディスクアレイ間ネットワークを有するディスク装置を提供できる。また、ディスクからの読み出しおよびディスクへの書き込みを高速スルーポート化できるディスク装置を提供できる。また、高スルーポートを維持できるディスクアレイ間ネットワークを有するディスクアダプタとディスクアレイ間のスルーポートが高く低コストなディスク装置を提供できる。

【図面の簡単な説明】

【図1】 本発明の第1の実施の形態のディスク装置を示す図である。  
【図2】 本発明に用いるスイッチの構成を示す図である。  
【図3】 本発明に用いるスイッチの構成を示す図である。  
【図4】 本発明に用いるスイッチの動作を示す図である。  
【図5】 本発明に用いるスイッチの動作を示す図である。  
【図6】 本発明の第1の実施の形態に対して、ディスクドライブを増設する方法を示す図である。

【図7】 本発明の第2の実施の形態のディスク装置を示す図である。

【図8】 本発明の第3の実施の形態のディスク装置を示す図である。

【図9】 従来のディスク装置を示す図である。

【図10】 FC-ALEによる接続形態を説明する図である。

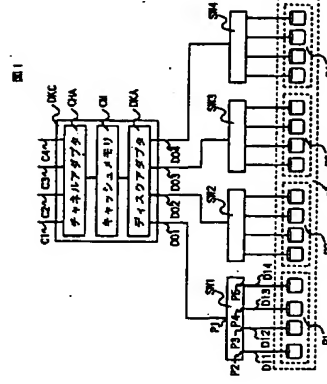
【図11】 FC-ALEによる接続形態を説明する図である。

【図12】 本発明の第4の実施の形態のディスク装置を示す図である。

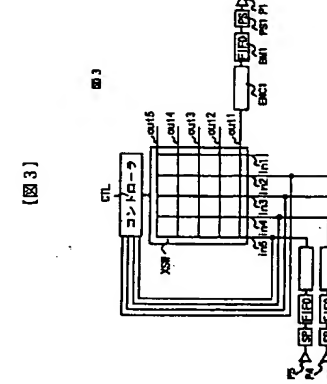
【符号の説明】

DKC, DKC1, DKC2 ディスクコントローラ  
CHA, CHA1, CHA2 チャネルアダプタ  
CM, CM1, CM2 キャッシュメモリ  
DKA, DKA1, DKA2 ディスクアダプタ  
DA, DA1~DA4, ディスクアレイ  
DK1~DK4 ディスクドライブ  
R1~R6 ディスクグループ  
C1~C4, D1~D4, D01~D05, D11~D14, D21~D24, D1a, D1b, D2a, D2b, D3a, D3b チャネル  
SW1~SW4 スイッチ  
P1~P5 スイッチポート  
XSW クロスバススイッチ  
CTL スイッチコントローラ  
In1~In5 クロスバススイッチ入力ポート  
Out1~Out5 クロスバススイッチ出力ポート  
SP1, SP2 シリアルパラレル変換装置  
PS1, PS2 パラレルシリアル変換装置  
BM1, BM2 バッファメモリ  
DEC1, DEC2 8B10B変換デコーダ  
ENC1, ENC2 8B10B変換エンコーダ  
T1, T2, T3,  
T4, T5, T6 フレームの時間  
Tx 送信機  
Rx 受信機  
NL NLポート  
FL FLポート

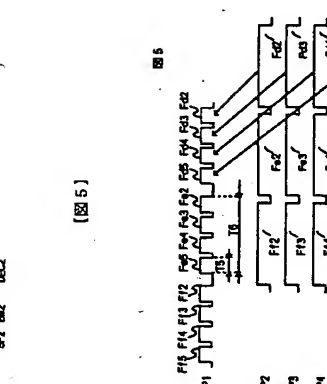
【図1】



【図2】



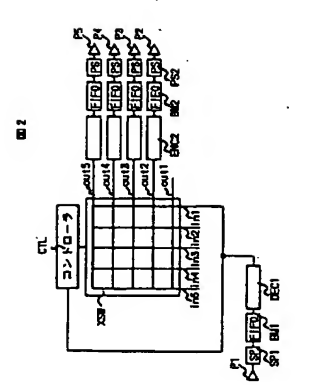
【図3】



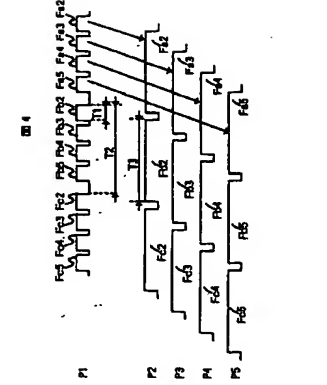
【図4】



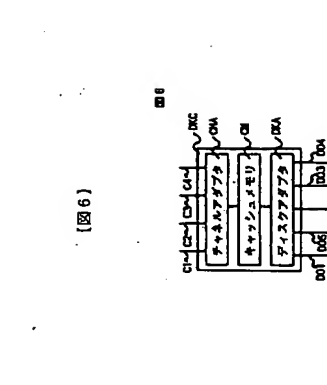
【図5】



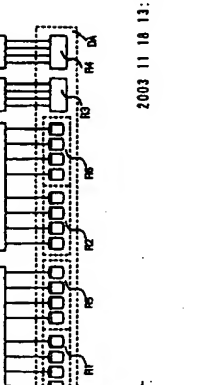
【図6】



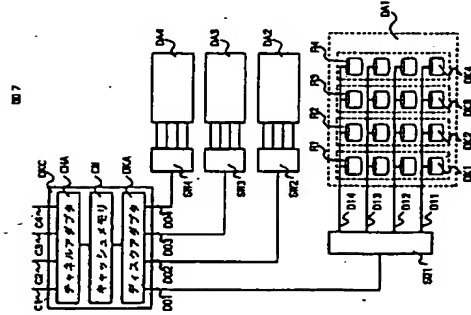
【図7】



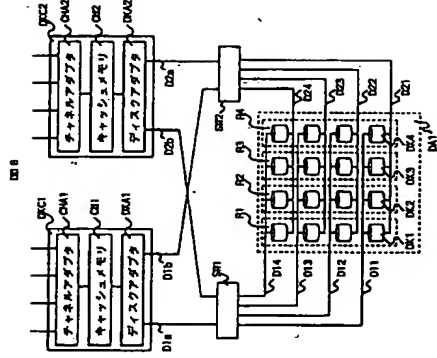
【図8】



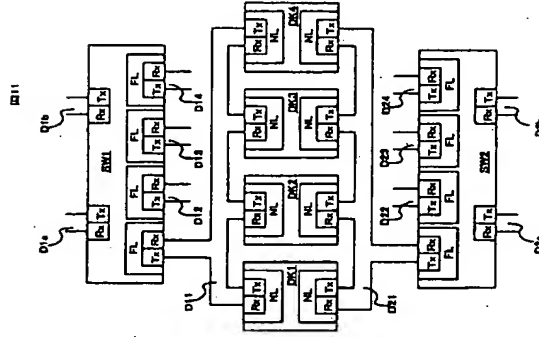
【図 7】



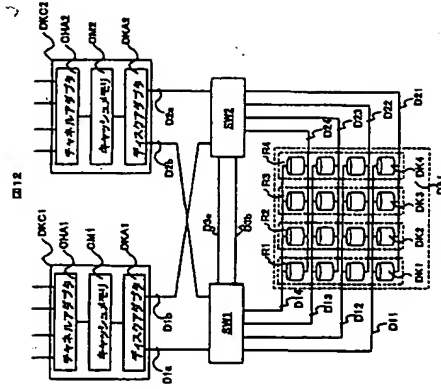
【図 8】



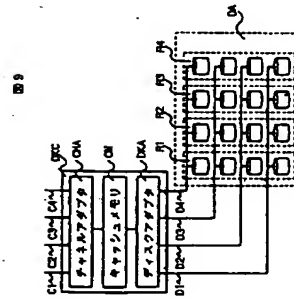
【図 11】



【図 12】



【図 9】



【図 10】

